

Introduction to Big Data Assignment-7

Name: Faizan Mulla

Roll No: 21F1003885

Problem Statement

1. Create a file with at least 1000 rows of data and upload it to a Google Cloud Storage (GCS) bucket.
2. Write a Kafka Producer application to:
 - Read data from the GCS bucket.
 - Process data in batches of 10 rows and send each batch to a Kafka topic.
 - Introduce a delay of 10 seconds between sending each batch.
3. Write a Spark Streaming Consumer application to:
 - Read data from the Kafka topic every 5 seconds.
 - Compute and emit the count of rows processed in the last 10 seconds.
4. Run the Producer and Consumer simultaneously, stopping after all 1000 rows are processed.

Solution

The solution involves the following key steps:

1. Setting up the environment locally and on GCP.
2. Implementing the Kafka Producer to read from GCS and send data to Kafka.
3. Writing the Spark Streaming Consumer to process Kafka messages.
4. Running and verifying both components simultaneously.

Implementation Details

Step 1: Environment Setup

Local Setup

1. Install required packages:

```
```bash
sudo apt update; sudo apt -y upgrade
sudo apt install -y python3.11-venv openjdk-11-jdk
```
```

```
```bash
python3 -m venv .venv
source .venv/bin/activate
pip install -U kafka google-cloud-storage pyspark pandas
```
```

2. Download and extract Kafka:

```
```bash
wget https://downloads.apache.org/kafka/3.8.0/kafka_2.13-3.8.0.tgz
tar -xzf kafka_2.13-3.8.0.tgz
rm kafka_2.13-3.8.0.tgz

cd kafka_2.13-3.8.0/
```
```

GCP Setup

1. Create a GCS Bucket

- Navigate to the Google Cloud Console.
- Create a new bucket named ``iitm-ibd-ga7``.
- Upload the 1000-row CSV file to this bucket.

2. Dataproc Cluster Creation

Created a Dataproc cluster with following specifications:

- chose the option: Create cluster on compute engine
- manager node: series → e2 // machine type → e2-standard-2 (2vCPU, 1 core, 8GB memory)
- reduce primary disk size from 500GB to something less like 50GB.
- exact same settings for worker nodes too.
- Region: **asia-south-1**
- in the customize cluster menu, **uncheck** the INTERNAL IP ONLY option.

Step 2: Kafka Producer Implementation

1. Producer Script (`kafka_producer.py`):

```
from kafka import KafkaProducer
import pandas as pd
import json
import time
from typing import List

class SalesDataProducer:
    def __init__(self, bootstrap_servers: str, topic_name: str):
        self.producer = KafkaProducer(
            bootstrap_servers=bootstrap_servers,
            value_serializer=lambda x: json.dumps(x).encode("utf-8"),
        )
        self.topic = topic_name

    def send_batch(self, records: List[dict]):
        for record in records:
            self.producer.send(self.topic, value=record)
        self.producer.flush()

    def process_file(self, file_path: str, batch_size: int = 10):
        df = pd.read_csv(file_path)
        records_processed = 0
        total_records = min(1000, len(df)) # Process maximum 1000 records

        while records_processed < total_records:
            batch_end = min(records_processed + batch_size, total_records)
            batch = df.iloc[records_processed:batch_end].to_dict("records")

            print(f"Sending batch of {len(batch)} records...")
            self.send_batch(batch)

            records_processed += len(batch)
            print(f"Total records processed: {records_processed}/{total_records}")

            if records_processed < total_records:
                print("Sleeping for 10 seconds...")
                time.sleep(10)

        self.producer.close()

if __name__ == "__main__":
    producer = SalesDataProducer(
        bootstrap_servers="localhost:9092", topic_name="sales_data"
    )
    producer.process_file("gs://iitm-ibd-ga7/sales_data.csv")
```

Step 3: Spark Consumer Implementation

1. Consumer Script (`spark_consumer.py`):

```
Tabnine | Edit | Test | Explain | Document | Ask
def create_spark_session():
    return (
        SparkSession.builder.appName("SalesDataConsumer")
        .config(
            "spark.jars.packages", "org.apache.spark:spark-sql-kafka-0-10_2.12:3.2.0"
        )
        .getOrCreate()
    )

Tabnine | Edit | Test | Explain | Document | Ask
def create_schema():
    return StructType(
        [
            StructField("timestamp", StringType(), True),
            StructField("user_id", IntegerType(), True),
            StructField("product_id", IntegerType(), True),
            StructField("quantity", IntegerType(), True),
            StructField("price", DoubleType(), True),
        ]
    )

Tabnine | Edit | Test | Explain | Document | Ask
def start_streaming():
    spark = create_spark_session()
    schema = create_schema()

    # Read from Kafka
    df = (
        spark.readStream.format("kafka")
        .option("kafka.bootstrap.servers", "localhost:9092")
        .option("subscribe", "sales_data")
        .option("startingOffsets", "earliest")
        .load()
    )

    # Parse JSON data and convert timestamp to timestamp type
    parsed_df = (
        df.selectExpr("CAST(value AS STRING)")
        .select(from_json(col("value"), schema).alias("data"))
        .select("data.*")
        .withColumn("timestamp", to_timestamp("timestamp", "yyyy-MM-dd HH:mm:ss"))
    )

    # Calculate counts over 10-second windows
    windowed_counts = (
        parsed_df.withWatermark("timestamp", "10 seconds")
        .groupBy(window("timestamp", "10 seconds", "5 seconds"))
        .agg(count("*").alias("record_count"))
    )

    # Start the streaming query
    query = (
        windowed_counts.writeStream.outputMode("update")
        .format("console")
        .option("truncate", "false")
        .trigger(processingTime="5 seconds")
        .start()
    )

    query.awaitTermination()

if __name__ == "__main__":
    start_streaming()
```

Step 4: Execution

In order to open a terminal, go to the running Dataproc Cluster and click on **VM instances**. Now, click on **SSH**, right next to the master node.

1. Start Kafka (Zookeeper server + Kafka server)

Open **Terminal 1**:

```
``bash
cd kafka_2.13-3.8.0/
source .venv/bin/activate
bin/zookeeper-server-start.sh config/zookeeper.properties
...
```

Open **Terminal 2**:

```
``bash
cd kafka_2.13-3.8.0/
source .venv/bin/activate

bin/kafka-server-start.sh config/server.properties
...
```

Create Kafka Topic: (*name=sales_data*)

```
``bash
bin/kafka-topics.sh --create --topic sales_data --bootstrap-server localhost:9092
--partitions 1 --replication-factor 1
...
```

2. Start producer service for `sales_data` topic

Open **Terminal 3**:

```
```bash
cd kafka_2.13-3.8.0/
source .venv/bin/activate
bin/kafka-console-producer.sh --topic sales_data --broker-list localhost:9092
```
```

3. Start **consumer service** for `sales_data` topic

Open **Terminal 4**:

```
```bash
cd kafka_2.13-3.8.0/
source .venv/bin/activate

bin/kafka-console-consumer.sh --topic sales_data --bootstrap-server
localhost:9092
```
```

4. Run **Spark Consumer**

Open **Terminal 5**:

- Upload the **spark_consumer.py** file using “Upload Files”
- Now run this file using:

```
```bash
source .venv/bin/activate
spark-submit --packages
org.apache.spark:spark-sql-kafka-0-10_2.12:3.5.0 spark_consumer.py
```
```

5. Run **Kafka Producer**

Open **Terminal 6**:

- Upload the *kafka_producer.py* file using “Upload Files”
- Now run this file using:

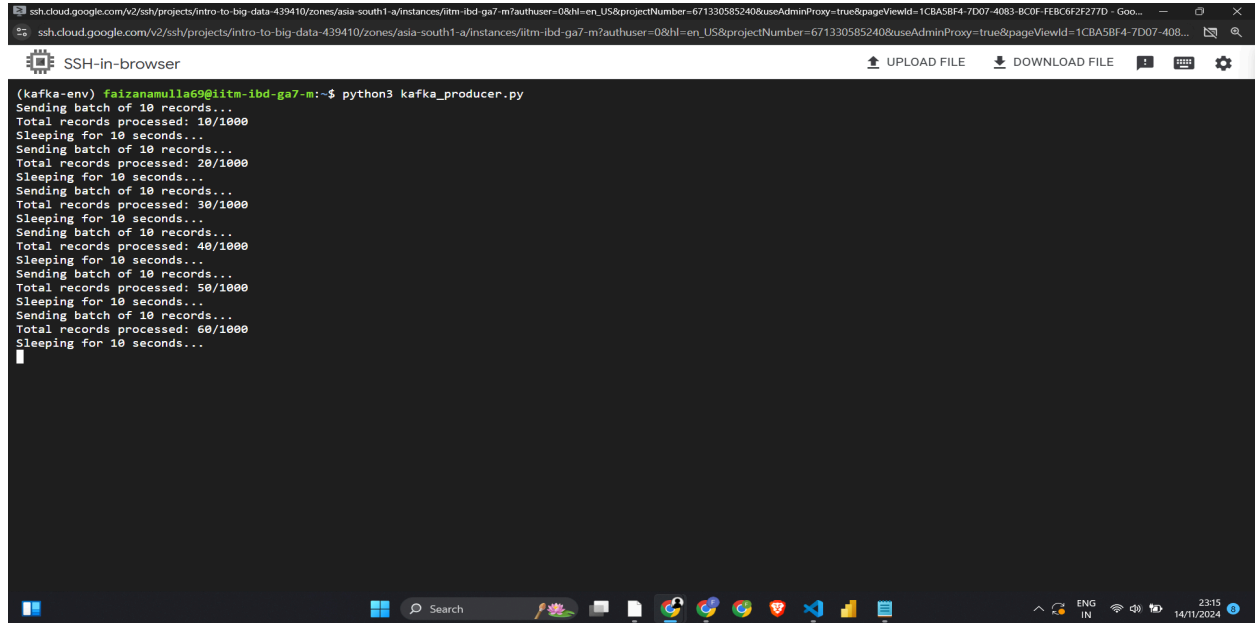
```
```bash
source .venv/bin/activate
python3 kafka_producer.py
```
```

Results

1. Kafka Producer successfully sent 1000 rows in batches of 10 with 10-second intervals.
2. Spark Consumer accurately computed and displayed row counts every 10 seconds.
3. The process stopped after all rows were processed.

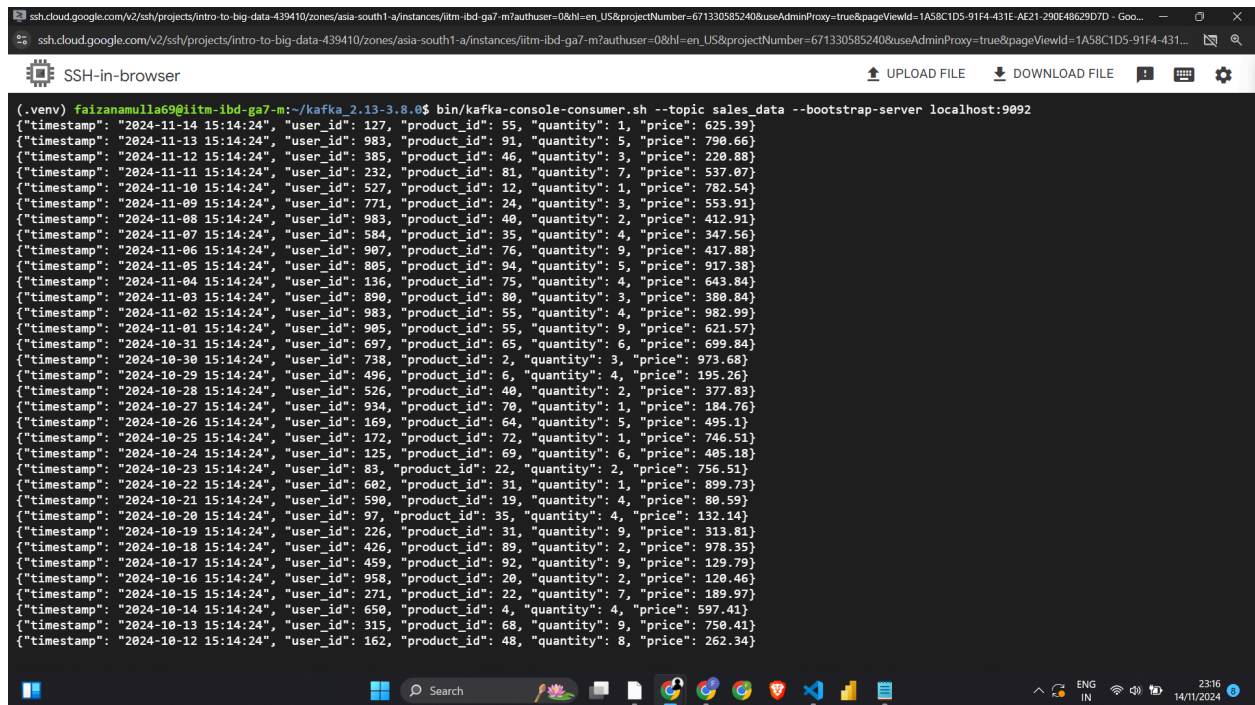
Relevant Screenshots

1. Kafka Producer Output:



```
(kafka-env) faizanamulla69@iitm-ibd-ga7-m:~$ python3 kafka_producer.py
Sending batch of 10 records...
Total records processed: 10/1000
Sleeping for 10 seconds...
Sending batch of 10 records...
Total records processed: 20/1000
Sleeping for 10 seconds...
Sending batch of 10 records...
Total records processed: 30/1000
Sleeping for 10 seconds...
Sending batch of 10 records...
Total records processed: 40/1000
Sleeping for 10 seconds...
Sending batch of 10 records...
Total records processed: 50/1000
Sleeping for 10 seconds...
Sending batch of 10 records...
Total records processed: 60/1000
Sleeping for 10 seconds...
```

2. Spark Consumer Output:



```
(.venv) faizanamulla69@iitm-ibd-ga7-m:~/kafka_2.13-3.8.0$ bin/kafka-console-consumer.sh --topic sales_data --bootstrap-server localhost:9092
{"timestamp": "2024-11-14 15:14:24", "user_id": 127, "product_id": 55, "quantity": 1, "price": 625.39}
{"timestamp": "2024-11-13 15:14:24", "user_id": 983, "product_id": 91, "quantity": 5, "price": 790.66}
{"timestamp": "2024-11-12 15:14:24", "user_id": 385, "product_id": 46, "quantity": 3, "price": 220.88}
{"timestamp": "2024-11-11 15:14:24", "user_id": 232, "product_id": 81, "quantity": 7, "price": 537.07}
{"timestamp": "2024-11-10 15:14:24", "user_id": 527, "product_id": 12, "quantity": 1, "price": 782.54}
{"timestamp": "2024-11-09 15:14:24", "user_id": 771, "product_id": 24, "quantity": 3, "price": 553.91}
{"timestamp": "2024-11-08 15:14:24", "user_id": 983, "product_id": 40, "quantity": 2, "price": 412.91}
{"timestamp": "2024-11-07 15:14:24", "user_id": 584, "product_id": 35, "quantity": 4, "price": 347.56}
{"timestamp": "2024-11-06 15:14:24", "user_id": 907, "product_id": 76, "quantity": 9, "price": 417.88}
{"timestamp": "2024-11-05 15:14:24", "user_id": 805, "product_id": 94, "quantity": 5, "price": 917.38}
{"timestamp": "2024-11-04 15:14:24", "user_id": 136, "product_id": 75, "quantity": 4, "price": 643.84}
{"timestamp": "2024-11-03 15:14:24", "user_id": 890, "product_id": 80, "quantity": 3, "price": 380.84}
{"timestamp": "2024-11-02 15:14:24", "user_id": 983, "product_id": 55, "quantity": 4, "price": 982.99}
{"timestamp": "2024-11-01 15:14:24", "user_id": 905, "product_id": 55, "quantity": 9, "price": 621.57}
{"timestamp": "2024-10-31 15:14:24", "user_id": 697, "product_id": 65, "quantity": 6, "price": 699.84}
{"timestamp": "2024-10-30 15:14:24", "user_id": 738, "product_id": 2, "quantity": 3, "price": 973.68}
{"timestamp": "2024-10-29 15:14:24", "user_id": 496, "product_id": 6, "quantity": 4, "price": 195.26}
{"timestamp": "2024-10-28 15:14:24", "user_id": 526, "product_id": 40, "quantity": 2, "price": 377.83}
{"timestamp": "2024-10-27 15:14:24", "user_id": 934, "product_id": 70, "quantity": 1, "price": 184.76}
{"timestamp": "2024-10-26 15:14:24", "user_id": 169, "product_id": 64, "quantity": 5, "price": 495.1}
{"timestamp": "2024-10-25 15:14:24", "user_id": 172, "product_id": 72, "quantity": 1, "price": 746.51}
{"timestamp": "2024-10-24 15:14:24", "user_id": 125, "product_id": 69, "quantity": 6, "price": 405.18}
{"timestamp": "2024-10-23 15:14:24", "user_id": 83, "product_id": 22, "quantity": 2, "price": 756.51}
{"timestamp": "2024-10-22 15:14:24", "user_id": 602, "product_id": 31, "quantity": 1, "price": 899.73}
{"timestamp": "2024-10-21 15:14:24", "user_id": 590, "product_id": 19, "quantity": 4, "price": 80.59}
{"timestamp": "2024-10-20 15:14:24", "user_id": 97, "product_id": 35, "quantity": 4, "price": 132.14}
{"timestamp": "2024-10-19 15:14:24", "user_id": 226, "product_id": 31, "quantity": 9, "price": 313.81}
{"timestamp": "2024-10-18 15:14:24", "user_id": 426, "product_id": 89, "quantity": 2, "price": 970.35}
{"timestamp": "2024-10-17 15:14:24", "user_id": 459, "product_id": 92, "quantity": 9, "price": 129.79}
{"timestamp": "2024-10-16 15:14:24", "user_id": 958, "product_id": 20, "quantity": 2, "price": 120.46}
{"timestamp": "2024-10-15 15:14:24", "user_id": 271, "product_id": 22, "quantity": 7, "price": 189.97}
{"timestamp": "2024-10-14 15:14:24", "user_id": 650, "product_id": 4, "quantity": 4, "price": 597.41}
{"timestamp": "2024-10-13 15:14:24", "user_id": 315, "product_id": 68, "quantity": 9, "price": 750.41}
{"timestamp": "2024-10-12 15:14:24", "user_id": 162, "product_id": 48, "quantity": 8, "price": 262.34}
```

3. GCS Bucket:

iitm-ibd-ga7

Location

asia-south1 (Mumbai)

Storage class

Standard

Public access

Not public

Protection

Soft Delete

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

OBSERVABILITY

INVENTORY REPORTS

OPERATIONS

Folder browser

iitm-ibd-ga7

Buckets > iitm-ibd-ga7

CREATE FOLDER

UPLOAD

TRANSFER DATA

OTHER SERVICES

Filter by name prefix only

Filter

Filter objects and folders

Show Live objects only

| <input type="checkbox"/> | Name | Size | Type | Created | Storage class | Last modified | |
|--------------------------|--------------------------------------|---------|----------|--------------------------|---------------|---------------|-----------------------------------|
| <input type="checkbox"/> | <div><div></div>sales_data.csv</div> | 35.8 KB | text/csv | Nov 14, 2024, 3:16:52 PM | Standard | Nov 14, 2024, | <div><div></div><div></div></div> |

4. Dataproc Cluster:

Cluster details

SUBMIT JOB

REFRESH

START

STOP

DELETE

VIEW LOGS

Consider using Auto Zone rather than selecting a zone manually. See <https://cloud.google.com/dataproc/docs/concepts/configuring-clusters/auto-zone>

Name

iitm-ibd-ga7

Cluster UUID

dc34531d-4441-4076-b783-7efaf9a6256b

Type

Dataproc Cluster

Status

Running

MONITORING

JOBS

VM INSTANCES

CONFIGURATION

WEB INTERFACES

Filter

Filter instances

| | Name | Role | Machine type |
|-----------------------------------|----------------------------------|--------|---------------|
| <div><div></div><div></div></div> | iitm-ibd-ga7-m | Master | e2-standard-2 |
| <div><div></div><div></div></div> | iitm-ibd-ga7-w-0 | Worker | e2-standard-2 |
| <div><div></div><div></div></div> | iitm-ibd-ga7-w-1 | Worker | e2-standard-2 |

EQUIVALENT REST